

REPORT DOCUMENTATION PAGE			Form Approved OMB NO. 0704-0188		
<p>The public reporting burden for this collection of information is estimated to average 1 hour per response, including the time for reviewing instructions, searching existing data sources, gathering and maintaining the data needed, and completing and reviewing the collection of information. Send comments regarding this burden estimate or any other aspect of this collection of information, including suggestions for reducing this burden, to Washington Headquarters Services, Directorate for Information Operations and Reports, 1215 Jefferson Davis Highway, Suite 1204, Arlington VA, 22202-4302. Respondents should be aware that notwithstanding any other provision of law, no person shall be subject to any penalty for failing to comply with a collection of information if it does not display a currently valid OMB control number.</p> <p>PLEASE DO NOT RETURN YOUR FORM TO THE ABOVE ADDRESS.</p>					
1. REPORT DATE (DD-MM-YYYY) 07-12-2015		2. REPORT TYPE Final Report		3. DATES COVERED (From - To) 1-Oct-2013 - 30-Jun-2014	
4. TITLE AND SUBTITLE Final Report: Socio-metrics: Identifying Invisible Deviant Adversaries			5a. CONTRACT NUMBER W911NF-13-1-0397		
			5b. GRANT NUMBER		
			5c. PROGRAM ELEMENT NUMBER 611102		
6. AUTHORS Gail-Joon Ahn			5d. PROJECT NUMBER		
			5e. TASK NUMBER		
			5f. WORK UNIT NUMBER		
7. PERFORMING ORGANIZATION NAMES AND ADDRESSES Arizona State University ORSPA P.O. Box 876011 Tempe, AZ 85287 -6011			8. PERFORMING ORGANIZATION REPORT NUMBER		
9. SPONSORING/MONITORING AGENCY NAME(S) AND ADDRESS (ES) U.S. Army Research Office P.O. Box 12211 Research Triangle Park, NC 27709-2211			10. SPONSOR/MONITOR'S ACRONYM(S) ARO		
			11. SPONSOR/MONITOR'S REPORT NUMBER(S) 64115-CS-II.1		
12. DISTRIBUTION AVAILABILITY STATEMENT Approved for Public Release; Distribution Unlimited					
13. SUPPLEMENTARY NOTES The views, opinions and/or findings contained in this report are those of the author(s) and should not be construed as an official Department of the Army position, policy or decision, unless so designated by other documentation.					
14. ABSTRACT In recent times, with the increasing growth in popularity of online social networks (OSNs) and Internet discussion forums, cybercriminals have found new ways to communicate and collaborate with each other in order to carry out cyber-attacks. Adversaries actively use Internet forums to form underground hacking communities where they exchange information on creating malicious programs and engage in the trade of malicious goods and services. Identifying the influential members of these underground communities who are behind the creation and distribution of tools used in cyber attacks would greatly help law enforcement agencies in controlling cybercrime. Manually					
15. SUBJECT TERMS Final Report for Seed Project					
16. SECURITY CLASSIFICATION OF:			17. LIMITATION OF ABSTRACT	15. NUMBER OF PAGES	19a. NAME OF RESPONSIBLE PERSON
a. REPORT UU	b. ABSTRACT UU	c. THIS PAGE UU			Gail-Joon Ahn
					19b. TELEPHONE NUMBER 480-/96-5900

Report Title

Final Report: Socio-metrics: Identifying Invisible Deviant Adversaries

ABSTRACT

In recent times, with the increasing growth in popularity of online social networks (OSNs) and Internet discussion forums, cybercriminals have found new ways to communicate and collaborate with each other in order to carry out cyber-attacks. Adversaries actively use Internet forums to form underground hacking communities where they exchange information on creating malicious programs and engage in the trade of malicious goods and services. Identifying the influential members of these underground communities who are behind the creation and distribution of tools used in cyber-attacks would greatly help law enforcement agencies in controlling cybercrime. Manually analyzing real-world data on hacking groups is tedious and requires enormous time and effort. For this seed project, we focus on Socia!SEAL, a tool which makes use of social network analysis techniques to reduce the manual effort required in identifying influential adversaries and visualizing the underlying social structure of underground hacking communities, that will eventually help identify links between attack attributions and influential adversaries in the next phase of this project.

Enter List of papers submitted or published that acknowledge ARO support from the start of the project to the date of this printing. List the papers, including journal references, in the following categories:

(a) Papers published in peer-reviewed journals (N/A for none)

Received

Paper

TOTAL:

Number of Papers published in peer-reviewed journals:

(b) Papers published in non-peer-reviewed journals (N/A for none)

Received

Paper

TOTAL:

Number of Papers published in non peer-reviewed journals:

(c) Presentations

Number of Presentations: 0.00

Non Peer-Reviewed Conference Proceeding publications (other than abstracts):

Received Paper

TOTAL:

Number of Non Peer-Reviewed Conference Proceeding publications (other than abstracts):

Peer-Reviewed Conference Proceeding publications (other than abstracts):

Received Paper

TOTAL:

Number of Peer-Reviewed Conference Proceeding publications (other than abstracts):

(d) Manuscripts

Received Paper

TOTAL:

Number of Manuscripts:

Books

Received Book

TOTAL:

Received Book Chapter

TOTAL:

Patents Submitted

Patents Awarded

Awards

Graduate Students

<u>NAME</u>	<u>PERCENT_SUPPORTED</u>
FTE Equivalent:	
Total Number:	

Names of Post Doctorates

<u>NAME</u>	<u>PERCENT_SUPPORTED</u>
FTE Equivalent:	
Total Number:	

Names of Faculty Supported

NAME

PERCENT SUPPORTED

FTE Equivalent:

Total Number:

Names of Under Graduate students supported

NAME

PERCENT SUPPORTED

FTE Equivalent:

Total Number:

Student Metrics

This section only applies to graduating undergraduates supported by this agreement in this reporting period

The number of undergraduates funded by this agreement who graduated during this period:

The number of undergraduates funded by this agreement who graduated during this period with a degree in science, mathematics, engineering, or technology fields:.....

The number of undergraduates funded by your agreement who graduated during this period and will continue to pursue a graduate or Ph.D. degree in science, mathematics, engineering, or technology fields:.....

Number of graduating undergraduates who achieved a 3.5 GPA to 4.0 (4.0 max scale):.....

Number of graduating undergraduates funded by a DoD funded Center of Excellence grant for Education, Research and Engineering:.....

The number of undergraduates funded by your agreement who graduated during this period and intend to work for the Department of Defense

The number of undergraduates funded by your agreement who graduated during this period and will receive scholarships or fellowships for further studies in science, mathematics, engineering or technology fields:.....

Names of Personnel receiving masters degrees

NAME

Total Number:

Names of personnel receiving PHDs

NAME

Total Number:

Names of other research staff

NAME

PERCENT SUPPORTED

FTE Equivalent:

Total Number:

Sub Contractors (DD882)

Inventions (DD882)

Scientific Progress

The observations made and results obtained through the analysis of the LiveJournal dataset using SocialSEAL. See attachment.

Technology Transfer

N/A

PROJECT REPORT

**Project Title: Socio-metrics: Identifying Invisible
Deviant Adversaries**

Project ID: W911NF-13-1-0397

**Project Duration: 11/ 01/2013 – 6/30/2014
(8 months, Seed Project)**

Sponsor: Army Reserch Office

Prepared by Prof. Gail-Joon Ahn, Arizona State University.

Executive Summary

In recent times, with the increasing growth in popularity of online social networks (OSNs) and Internet discussion forums, cybercriminals have found new ways to communicate and collaborate with each other in order to carry out cyber-attacks. Adversaries actively use Internet forums to form underground hacking communities where they exchange information on creating malicious programs and engage in the trade of malicious goods and services. Identifying the influential members of these underground communities who are behind the creation and distribution of tools used in cyber-attacks would greatly help law enforcement agencies in controlling cybercrime. Manually analyzing real-world data on hacking groups is tedious and requires enormous time and effort. For this seed project, we focus on *SocialSEAL*, a tool which makes use of social network analysis techniques to reduce the manual effort required in identifying influential adversaries and visualizing the underlying social structure of underground hacking communities, that will eventually help identify links between attack attributions and influential adversaries in the next phase of this project.

1 Introduction

Online Social Networks (OSNs) and Internet discussion forums have become very popular over the last decade and are now used by millions of people to communicate and collaborate with each other. Since communication over the Internet is largely anonymous, it has become a popular medium among adversaries, who have been found to be actively leveraging Internet forums and OSNs in recent years to share knowledge on developing malicious programs and engage in the trade of malicious goods and services such as stolen credit card numbers, forged identification documents, malicious bots and so on [1, 2].

Trade through the Internet is a far more profitable option for adversaries compared to offline media due to the greater reach provided by online advertisements. Furthermore, the use of OSNs and Internet forums has allowed adversaries from around the world to form underground groups and collaborate each other in conducting cyber-attacks at an unprecedented scale [3]. Monitoring OSNs and Internet forums is thus imperative to understand and investigate cybercrime.

Previous research has shown that only a small subset of the individuals in a hacking community are actively engaged in developing new hacking techniques and tools for exploiting vulnerabilities in software products while the vast majority comprises individuals commonly known as "Script Kiddies" [4] who participate in cyber-attacks by just purchasing and executing the available hacking tools. Identifying these influential members of the underground hacking groups who are behind the creation and distribution of new exploits is thus necessary in order to thwart cybercrime.

One way of identifying influential adversaries in hacking groups on OSNs and discussion forums is to manually analyze the posts, comments and social circles of users of these websites. However, such manual analysis is quite tedious and can be very expensive and infeasible with

large data sets. Hence, there is a need to build systems which can reduce the manual effort involved in analyzing large hacking communities and for this purpose we have designed and developed *SocialSEAL*, a tool which makes use of social network analysis techniques to identify influential adversaries and visualize the social structure of underground hacking communities.

In Section 2 of this report, we explain the framework and implementation of *SocialSEAL* and in Section 3, we present the results of using *SocialSEAL* on data collected from real-world OSNs. Section 4 concludes this report including several areas where this work can be further improved in future projects.

2 *SocialSEAL*: Framework, Design and Implementation

SocialSEAL makes use of social network analysis techniques which are based on the *SocialImpact* framework discussed in [5]. This tool is designed to provide end-users both pre-computed statistics computed from the underlying data and statistics generated dynamically from the data through external user inputs. A detailed explanation on these analyses is elaborated in subsequent sections (Sections 2.1 and 2.2).

2.1 Basic Model and Pre-Computed Statistics

An OSN can be represented by six fundamental entities and five basic types of unidirectional relationships between them. Users are those who have profiles in the network and have the rights to join groups, post articles and give comments to others. Groups are those to which users can belong. In an OSN, groups are mainly formed based on common interests. Articles are posted by users who want to share them with the society. In an OSN, articles might introduce the latest technologies, analyze recent vulnerabilities, call for participation of network attacks and trade newly developed and deployed botnets. Comments are the subsequent posts to articles. Posts are the union of articles and comments. Strings are the elementary components of articles and comments. Strings are not necessarily meaningful words. They could be names, URLs and underground slangs. A user has a relationship *authorOf* with each post she/he authored. A user has a relationship *followerOf* with each user she/he follows. A user has a relationship *memberOf* with each group she/he joins. An article has a relationship *hostOf* with each comment it receives. A post has a relationship *containerOf* with each string it consists of. The following formal description summarizes the above-mentioned entities and relationships.

Definition 1 (Online Social Dynamics). An OSN is modeled with the following components:

- U, G, A, C, P, S denote a set of users, user groups, articles, comments, posts ($P = A \cup C$) and strings, respectively;
- $UP = \{(u, p) \mid u \in U, p \in P \text{ and } u \text{ has an } \textit{authorOf} \text{ relationship with } p\}$ is a one-to-many user-to-post relation denoting a user and her posts;
- $FL = \{(u, y) \mid u \in U, y \in U \text{ and } u \text{ has a } \textit{followerOf} \text{ relationship with } y\}$ is a many-to-many user-to-user follow relation;
- $MB = \{(u, g) \mid u \in U, g \in G \text{ and } u \text{ has a } \textit{memberOf} \text{ relationship with } g\}$ is a many-to-many user-to-group membership relation;

- $AC = \{(a,c) \mid a \in A, c \in C \text{ and } a \text{ has a hostOf relationship with } c\}$ is a one-to-many article-to-comment relation denoting an article and its following comments; and
- $PS = \{(p,s) \mid p \in P, s \in S \text{ and } p \text{ has a containerOf relationship with } s\}$ is a many-to-many post-to-string relation.

In order to help examiners identify adversarial behaviors in OUSDs, we also need to address the following critical issues related to evidence mining in underground society: How can we identify unintended behaviors among the crowd of social users? Given the additional evidence acquired from other sources, how can we correlate them with social dynamics in the identified hidden networks? How can we measure the evolution in such an underground community? To answer these questions, we articulate several principles that the measures for underground social dynamics analysis should follow:

Principle 1 The measures should support identifications of interesting adversaries and groups based on both their social relationships and online conversations.

Principle 2 The measures should be able to take external evidence into account and support interactions with security analysts.

Principle 3 The measures should support temporal analysis for the better understanding of the evolution in adversarial society.

Based on the abovementioned principles, the pre-computed statistics provided by *SocialSEAL* are User Influence, User Activeness, Group Influence, Group Activeness and Risk Index. Although these statistics are pre-computed and displayed when the tool is launched, the end-user has controls in adjusting parameters involved in the computation of these statistics based on the end-user's needs.

2.1.1 User & Group Influence

User Influence (UI) is calculated as the weighted sum of four user attributes: the total length of all of the user's posts (LP), the number of comments received by the user for his/her posts (CR), the number of external links in the user's posts (EL) and the number of followers the user has in his/her social network (FLR). The attribute LP provides a measure of how knowledgeable a user might be. It is usually observed that people with great knowledge in a particular field tend to author posts which are more elaborate compared to posts made by people with limited knowledge in the same field [6]. Likewise, the attribute EL measures the novelty of a user's posts. Novelty is less likely to contain many external links [6]. The attributes CR and FLR measure the popularity of the user in his/her social sphere. If $Z = \langle LP, EL, CR, FLR \rangle$ is defined as a vector having these four user attributes as its components and W is a column vector containing the corresponding weights for each of these components, then the User Influence (UI) is calculated as

$$UI = W^T Z$$

Group Influence (GI) is just calculated as the average User Influence (UI) computed over all the n users belonging to a hacking community.

$$GI = \left(\frac{1}{n}\right) * \sum UI$$

The User & Group Influence scores are normalized to lie between 0 and 1.

2.1.2 User & Group Activeness

User Activeness (UA) is calculated as the weighted sum of three user attributes: the number of posts and comments made by the user (PC), the number of people the user follows (FLW) and the number of groups the user is a member of (GM). If $Y = \langle PC, FLW, GM \rangle$ is defined as a vector with these attributes as its components and W is a column vector containing the corresponding weights for each of these components, then User Activeness (UA) is defined as:

$$UA = W^T Y$$

Group Activeness (GA) is just calculated as the average User Activeness (UA) computed over all the n users belonging to a hacking community.

$$GA = \left(\frac{1}{n}\right) * \sum UA$$

The User & Group Activeness scores are also normalized to lie between 0 and 1.

2.1.3 User & Group Risk Index

User Risk Index (URI) measures how risky a particular user is by making use of two attributes calculated from the user's posts: the number of unique risky terms used by the user (NRT) and the frequency of these risky terms (FRT). We define a risky term (RT) as any term which appears in attack attributions containing cybersecurity terms and names of popular malicious programs. The attribute NRT measures the number of topics related to cybersecurity that a user has interest in and the attribute FRT measures the extent to which the user is interested in these topics. If $R = \langle NRT, FRT \rangle$ is a vector having these two attributes as its components and W is a column vector containing the corresponding weights for these components. User Risk Index (URI) is computed as

$$URI = W^T R$$

Group Risk Index (GRI) is calculated as the average User Risk Index (URI) computed over all the n users belonging to a hacking community. Also, The Risk Index scores are normalized to lie between 0 and 1.

$$GRI = \left(\frac{1}{n}\right) * \sum URI$$

2.2 Dynamically Generated Statistics

The dynamically generated statistics provided by *SocialSEAL* are User Relevance and Group Relevance. *SocialSEAL* provides end-users the ability to query the underlying data on hacking communities collected from real-world OSNs and discussion forums. Whenever a query is issued

by the end-user, files from the data, which are relevant to the query, are displayed along with a list of the users and groups in the data ranked based on their relevance to the query issued.

2.2.1 User & Group Relevance

User Relevance (UR) measures how relevant a user is to a query issued by the end-user in real-time. This is calculated based on results generated by a search engine which uses Cosine Similarity to provide all user posts relevant to the query [7]. Based on the results provided by the search engine, it is possible to determine two attributes: the number of times the user has made use of the query terms in his/her posts (NQTU) and the number of times the query terms have appeared in other user/group posts which a user has subscribed to (NQTO). If $Q = \langle NQTU, NQTO \rangle$ is a vector having these two attributes as its components and W is a column vector containing the weights corresponding to these attributes. Therefore, the User Relevance (UR) is computed as

$$UR = W^T Q$$

The Group Relevance (GR) is just calculated as the average User Relevance (UR) computed over all the n users belonging to a hacking community. The User & Group Relevance scores are also normalized to lie between 0 and 1.

$$GR = \left(\frac{1}{n}\right) * \sum UR$$

2.3 System Architecture and Implementation

SocialSEAL has been designed to function as a web application operating under a three-tier architecture. Figure 1 shows the architecture of the *SocialSEAL* system.

Data Storage Layer The Data Storage Layer forms the base of the system and it consists of a graph database that stores information extracted from the data on hacking communities, which is collected from real-world OSNs and Internet forums. The data comprises HTML files containing profile and post information of the members of the hacking communities. These files are parsed using HTML parsers which are built for each website that is crawled for data. The information obtained by parsing these files is pre-processed in order to extract various useful attributes and meta-data, which are then stored in the graph database. Some of the attributes extracted include - username, location, number of followers in the user's social network, user's self-reported interests, and number of user posts.

In the current version of *SocialSEAL*, the database contents are static¹. We chose the Neo4j graph database in the database layer instead of relational database systems like MySQL since it has been found that graph databases work well on highly connected data (which is the case with OSNs). Moreover, the performance of graph algorithms is claimed to be much faster on graph databases compared to relational databases, which would need expensive join queries [8]. The

¹ In the next phase of this project, we plan to make it dynamic by allowing end-users to upload new data to the database

Data Storage Layer also has a search index generated over the files in the dataset using the Apache Lucene library. This index is used by *SocialSEAL*'s search engine which end-users could use to search for specific security-related terms in profiles and posts of the users and communities in the dataset crawled.

Computation Layer At the middle level of the system is the Computation Layer, which consists of an Apache Tomcat Server instance running the *SocialSEAL* web application. We also used the Java Vaadin framework, which is built upon the Google Web Toolkit (GWT) since it supports rapid application development, provides the ability to build professional UIs and also scales well. Compared to Java Web framework, it has been observed that Vaadin has shown the better performance [9].

The Computation Layer's functionality can be split across two modules:

1. *Intelligence Generator*: This module takes user-provided configuration values, if any, and computes the pre-computed statistics such as Influence, Activeness and Risk Index, which were explained in Section 2.1.
2. *Search Engine*: This module comprises a Search Engine using Apache Lucene, which uses the Vector Similarity measure to rank and return search results relevant to user-generated queries and also computes the dynamic statistics including Relevance, which was discussed in Section 2.2.

Data Visualization & Analysis Layer The Data Visualization & Analysis Layer forms the top level of the system and it supports the UI, data visualizations and end-user controls, which are used for analysis. This layer comprises the following modules:

1. *Data Visualization Viewer*: This module runs visualization scripts using the D3.js JavaScript library, which is a very well-known library for creating good-looking data visualizations. In addition, *SocialSEAL* currently provides two kinds of data visualizations:
 - i. The Force-Layout visualization [10], which is used to display the social network structure of the hacking groups in the dataset and for showing the User/Group Influence and Activeness rankings visually.
 - ii. The Word-Cloud visualization [11], which gives the end-user general idea of the kind of words used by hackers in their profile and posts.
2. *Analysis Control*: This module provides the end-user with controls for tweaking the parameters in the Influence, Activeness, Risk Index and Relevance score computations.
3. *Query Input*: This module takes queries from the end-user and passes them down to the Search Engine module in the Computation Layer.
4. *Results Viewer*: This module displays the results from the Influence, Activeness, Risk and Relevance Computations.

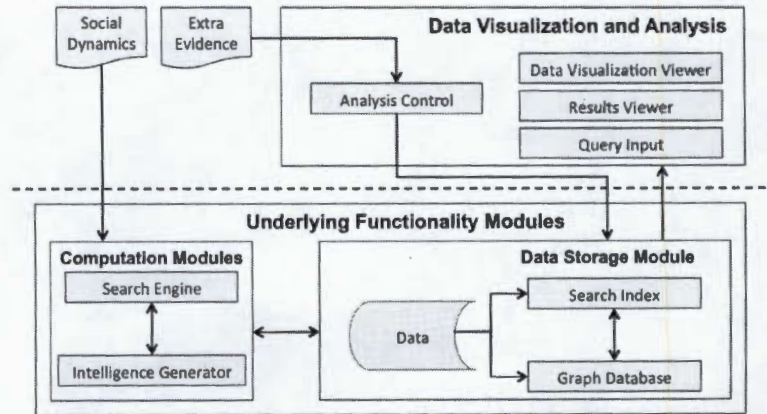


Figure 1. System Architecture of SocialSEAL

3 SocialSEAL in Action: Real-world Data Analysis

To evaluate *SocialSEAL*, we used a dataset crawled from a real-world OSN which contains the profiles and posts of some popular hacking communities and their members. It consists of more than 20,000 articles, 6,000 users, and 4,000 groups collected over 4 years. We analyzed the dataset using the various features provided by *SocialSEAL*. A description of the dataset used for the evaluation is provided in Section 3.1, followed by the pre-processing steps applied to extract useful attributes from the dataset in Section 3.2. The observations and results of our analysis of the dataset using SocialSEAL are described in Section 3.2.

3.1 Description of the dataset used

We collected a dataset containing the profiles and posts of five popular hacking communities and their members from LiveJournal, a popular Russian Social Network. Since there has been a spate of cyber-attacks from Russia and Eastern Europe in recent times [12], we felt that the possibility of finding adversaries would be much higher on Russian language forums and OSNs compared to more popular OSNs like Twitter or Facebook.

There are a total of 168 users in the dataset spread across the five hacking groups. Table 1 lists the hacking groups crawled and also shows the number of users in each group. For each user in the five targeted hacking groups, we crawled the user's profile page, pages containing posts authored by the user and pages containing the updates the user has received from his/her social circle. The profile and posts for each of the five hacking groups was also collected. The data collected is in the form of HTML files and it contains information generated by users between the years 2004 and 2008.

Hacking Group	Total Members	Influence Rank	Activeness Rank	Risk Index Rank
BH Crew	99	4	3	3
CUP SU	17	3	4	4
Damagelab	27	5	5	5
Mazafaka.ru	15	2	2	2
RU Hack	10	1	1	1

Table 1. Statistics on the hacking groups crawled

3.2 Data Pre-processing

Although a user's post can contain image and video data in addition to text, we only analyze the text of a user's post in the current version of *SocialSEAL*. Since the majority of users on LiveJournal are Russian, a lot of the user posts in the dataset had to be translated using the Microsoft Bing Translator service. Using data pre-processing scripts, we extracted several attributes from each user's profile and posts in order to compute the statistics listed under Section 2 of this report. The attributes extracted from the dataset are listed in Table 2.

Attributes collected from user profile	Attributes collected from user posts
1. Username and real name	1. Length of Posts
2. Location	2. Number of external links shared
3. Personal Website URL	3. Number of unique risky terms used
4. Interests	4. Total number of risky terms used
5. Number of followers	
6. Number of users followed	
7. Number of groups joined	
8. Total posts authored	
9. Total comments received	

Table 2. List of attributes extracted from the dataset

For each user in the dataset, we created a '*User*' node in the graph database and this node was populated with the attributes extracted from the dataset. A link between two '*User*' nodes in the database indicates that the users are related to each other. Similarly, we created '*Group*' nodes for each of the five groups being analyzed and these nodes were also populated with the attributes extracted from the dataset. A link between a '*User*' node and a '*Group*' node indicates that the user is a member of the group. Using link labels, we were able to specify the nature of the relationship. For example, '*followed by*', '*follows*' etc. are links between two '*User*' nodes and '*member of*', '*contributes to*' etc. are links between a '*User*' node and a '*Group*' node. Figure 2 shows the social graph constructed from the dataset using these concepts. The nodes with dark circles around them represent the '*Group*' nodes and the other nodes are user nodes. The yellow links represent '*User-User*' relationships and maroon links represent '*User-Group*' relationships. The nodes are color-coded to indicate which group they belong to.

We also constructed an index over the pages containing posts authored by each user and posts authored by each user's friends in his/her social circle. This index is used by the Search Engine module described earlier in Section 2.3, which can be used by the end-user to query the dataset for pages relevant to the terms specified in the user-provided query.

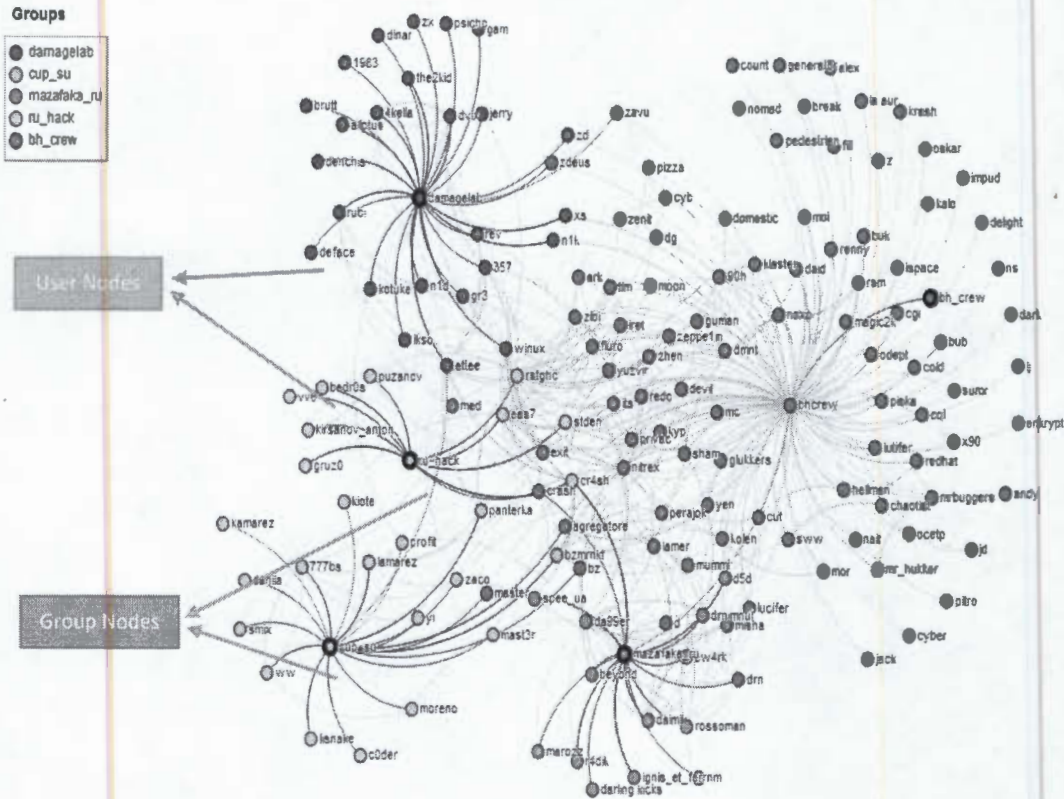


Figure 2. Visualization of the social network in the dataset

3.3 Analysis Results

The observations made and results obtained through the analysis of the LiveJournal dataset using *SocialSEAL* are described in the following sections.

3.3.1 Influence & Activeness Analysis

SocialSEAL can display the users and groups in the dataset ranked by their Influence and Activeness scores. We had discussed how these measures are computed under Section 2.1 earlier. The Influence and Activeness ranks of the five hacking groups in the dataset have been shown under Table 1. Similarly, the top 5 influential and active users in the dataset along with their respective scores are listed in Table 3.

Top 5 Influential Users	Influence Score	Top 5 Active Users	Activeness Score
1. eas7	1.0	1. ark	1.0
2. kalo	0.82	2. gam	0.595
3. ark	0.761	3. jd	0.575
4. fill	0.732	4. kolen	0.556
5. pitro	0.635	5. med	0.479

Table 3. Top 5 Influential & Active Users and their scores

From Table 3, it can also be seen that there is little correlation between the Top 5 influential and active users. This gives the idea that individuals who are very active in the hacking groups might not be influential people. However, the Influence and Activeness ranks of the hacking groups seem to be very similar, as can be noticed from Table 1.

Figures 3 and 4 show the nodes in the dataset ranked according to their Influence & Activeness scores. These figures also show the controls provided to the end user which can be used to tweak the weights of the parameters used in the Influence and Activeness computation. The variation in the 'User'/'Group' node size gives an intuitive idea of the Influence/Activeness rank of a user/group.

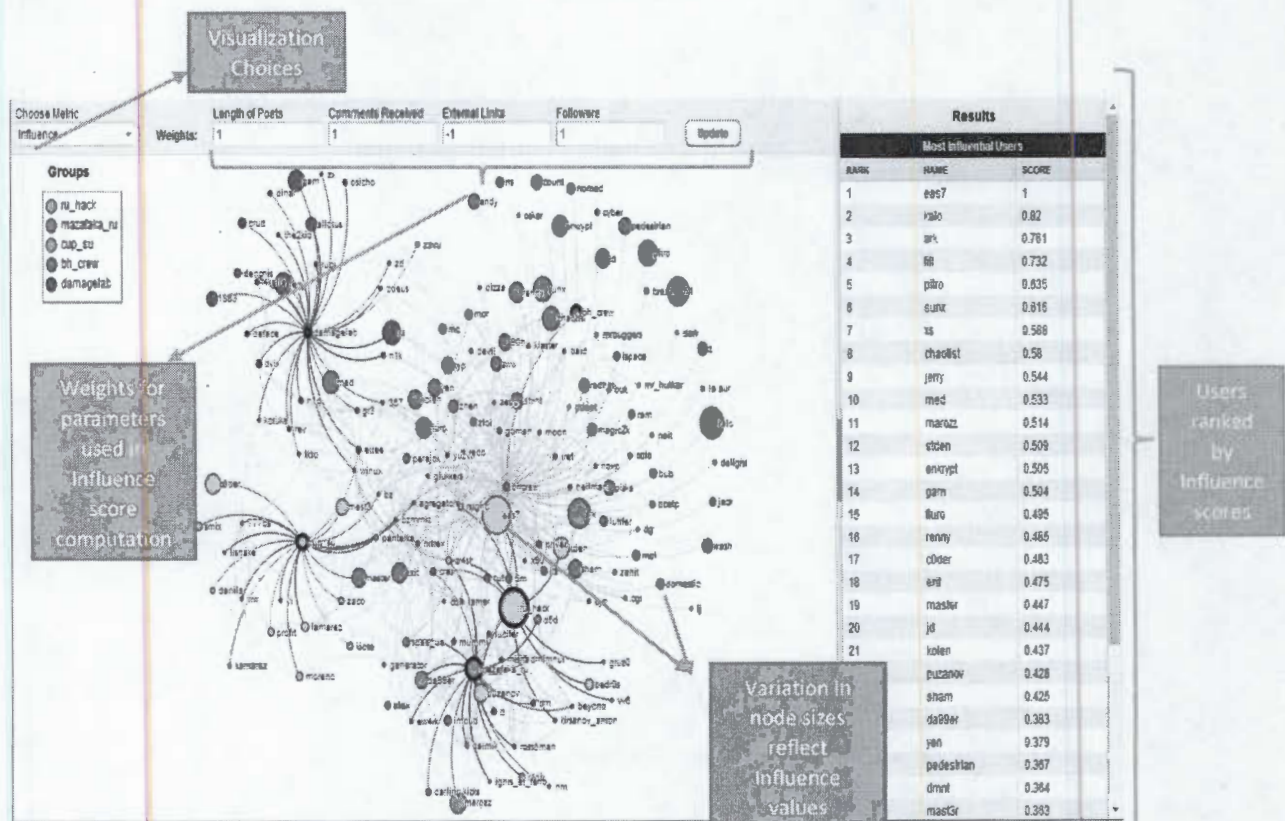


Figure 3. Nodes ranked by Influence

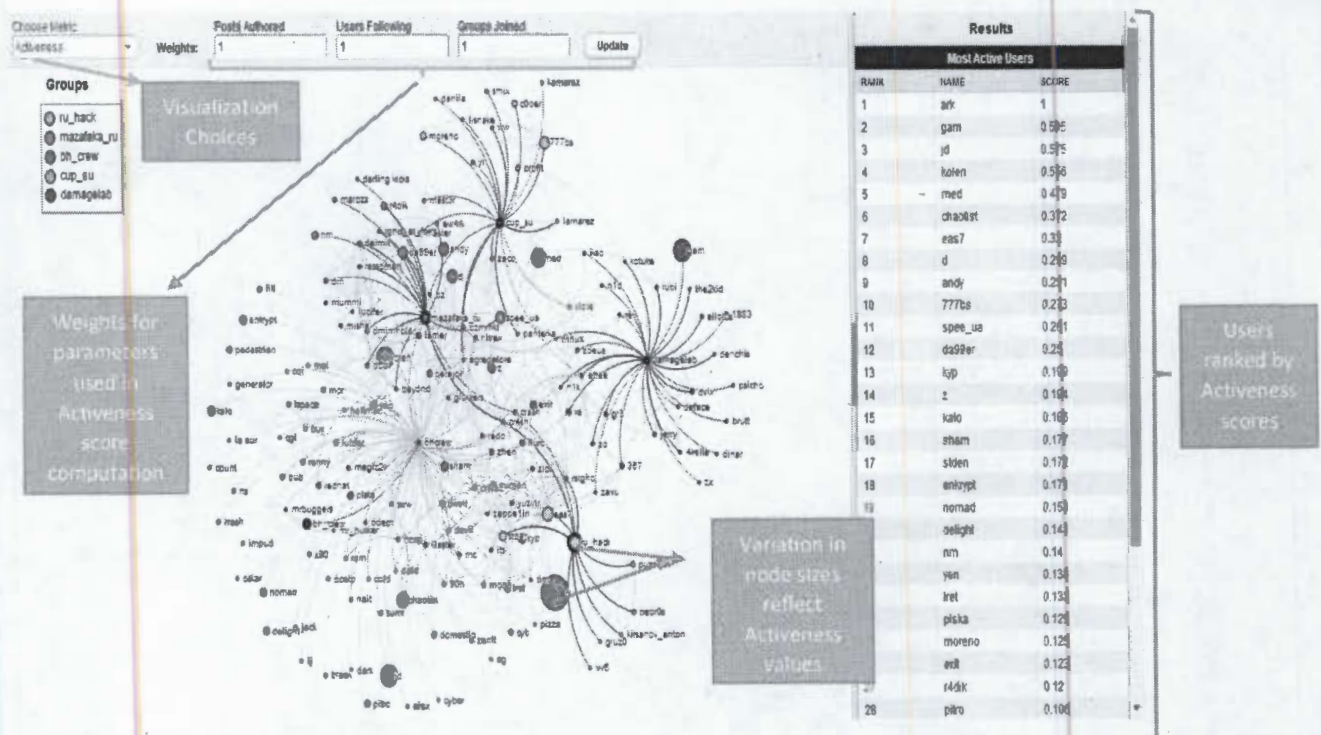


Figure 4. Nodes ranked by Activeness

3.3.2 Relevance Analysis

Using the querying feature in *SocialSEAL*, the end-user can query the built-in search engine for user pages containing posts similar to the query being presented. In addition to the search results, the search engine also returns the users and groups in the dataset ranked according to their relevance to the query. Since most of the content in the dataset is in Russian, the query input module issues the original query along with its Russian equivalent to the search engine.

Table 4 shows the top 5 relevant users and the groups they belong to for five sample queries issued to the search engine: botnet, virus, vulnerability, threat and rootkit. The numbers in brackets next to each query represent the number of search hits for the query. On comparing the results for the five sample queries, we observed that a few users appear among the top 5 relevant users for multiple queries. This shows who could be potentially influential adversaries in the dataset.

Botnet (23)	Virus (101)	Vulnerability (140)	Threat (47)	Rootkit (58)
1. exit	1. exit	1. eas7	1. sham	1. eas7
2. cr4sh	2. stden	2. sham	2. tim	2. cr4sh
3. crash	3. tim	3. zloi	3. ark	3. crash
4. dmnt	4. c0der	4. nait	4. exit	4. sww
5. smix	5. cr4sh	5. ark	5. its	5. ark

Table 4. Top 5 relevant users for five sample queries. Numbers in brackets represent total search hits.

Figure 5 shows the nodes in the dataset ranked according to their relevance to the query “botnet”. Similar to the interface used in the Influence/Activeness computation as illustrated in Figures 3 and 4, it enables the end user to adjust the weights of the parameters used in the Relevance computation in this interface as well. On the sidebar, both the search results and the users/groups in the dataset ranked by their Relevance scores are shown. Also, the variation in node sizes is used to give an intuitive idea of User/Group Relevance.

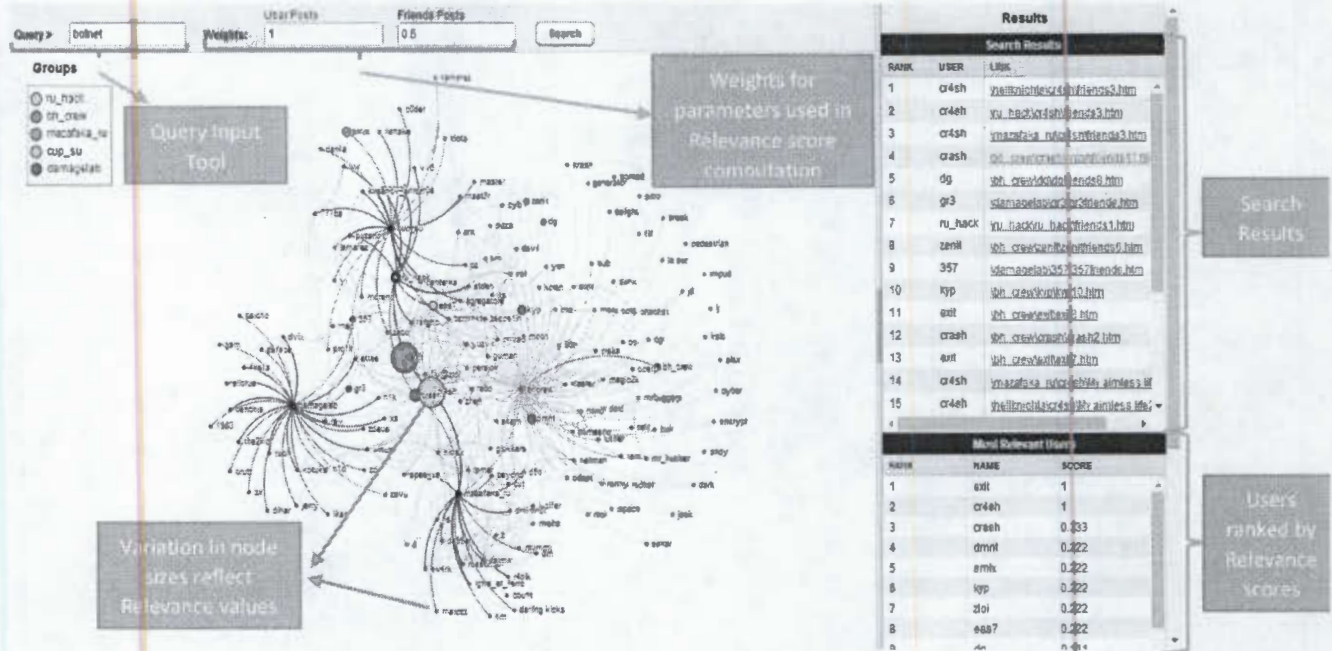


Figure 5. Nodes ranked by relevance to the query "botnet"

3.3.3 Risk Analysis

Based on the presence of keywords related to cybersecurity in the content posted by users, *SocialSEAL* generates a score for each user which measures how risky the user could be. A dictionary consisting of words related to cybersecurity was compiled from various sources on the Internet [13,14] and this dictionary was used by *SocialSEAL* to generate the risk index score, which was explained in Section 2.1. Table 5 shows the top 5 users in the dataset ranked according to their risk scores. Similarly, Table 3 shows the risk index ranks for the hacking groups in the dataset. It also allows the end user to tweak the weights of the parameters involved in the computation of the risk scores. Figure 6 shows the nodes in the dataset ranked according to their risk index scores.

Top 5 Risky Users	Risk Index Score
1. eas7	1.0
2. stden	0.634
3. puzanov	0.549
4. exit	0.544
5. kalo	0.439

Table 5. Top 5 Risky users and their scores

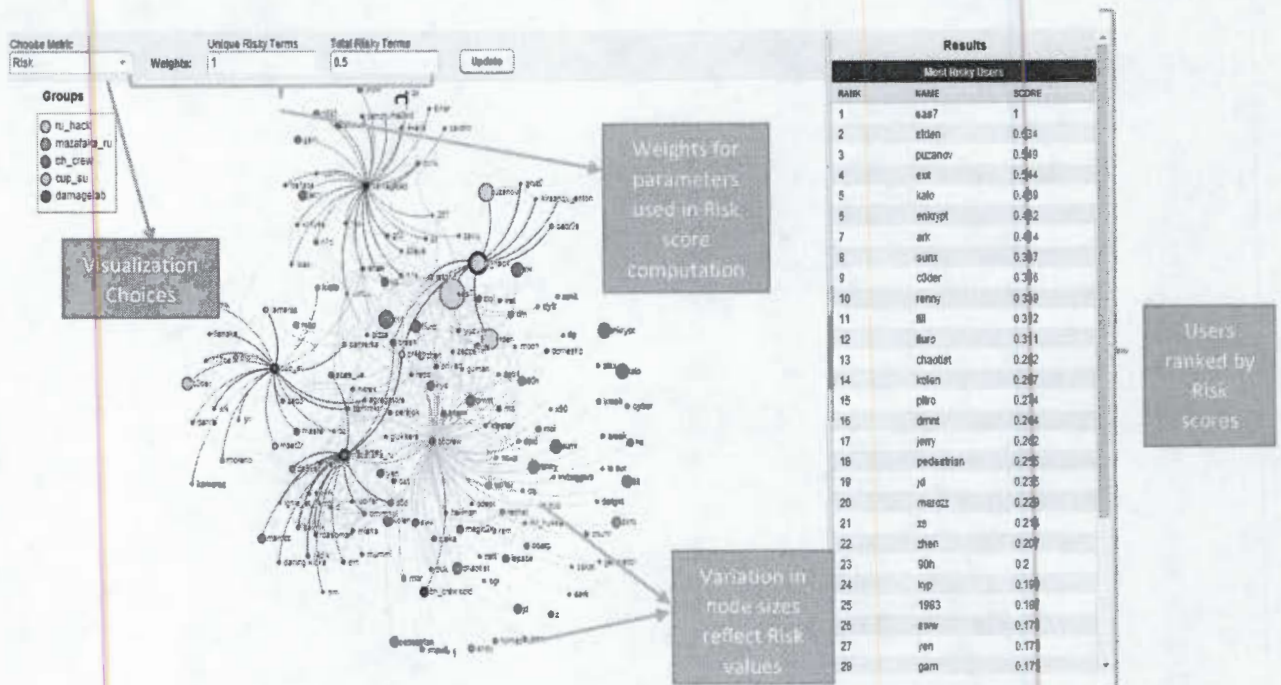
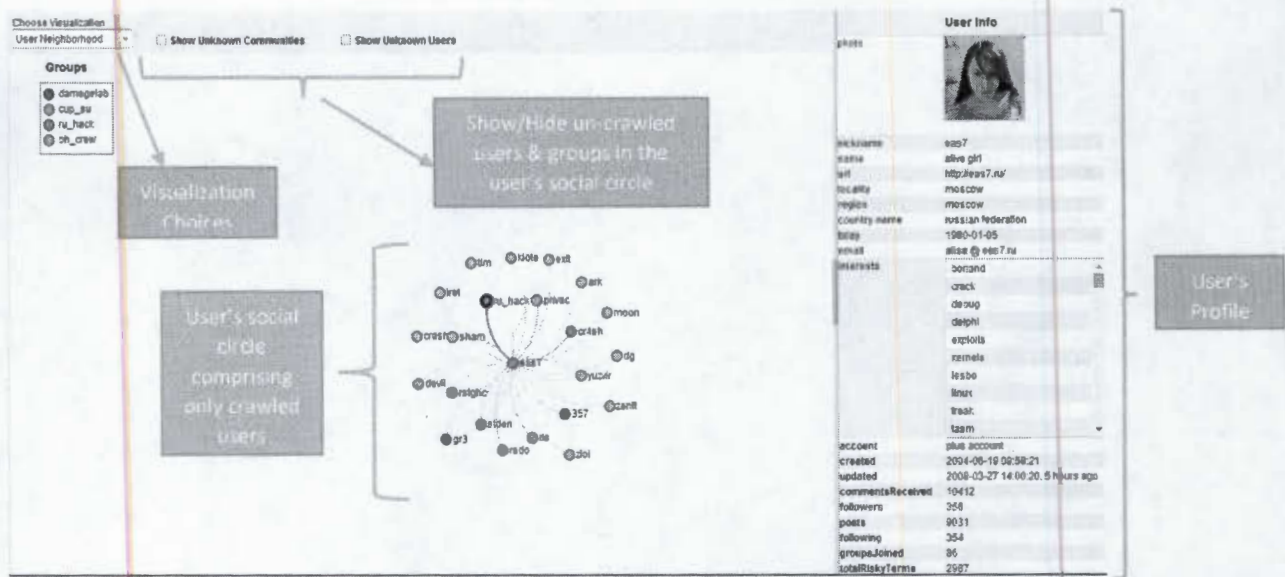


Figure 6. Nodes ranked according to Risk

On clicking any User or Group node in the graph, *SocialSEAL* takes the end-user to a page which shows just the user's social circle and information from the user's profile such as user's photo, location, interests etc. Since we have only crawled LiveJournal users belonging to the five targeting hacking communities, there could be other users/communities in the selected user's social circle who have not been crawled. Hence, there are two 'Show/Hide' check boxes which show/hide un-crawled users and communities. Figure 7 shows the profile of the most risky user in the dataset, who goes by the name **eas7**. The Word-Cloud generated by using the top 1,000 terms from this user's posts filtered based on the dictionary used in computing the Risk Index score is shown in Figure 8. In addition, we were able to check other security-related terms such as zombie, crack, hack, rootkit, spam, exploit, attack, botnet, and shellcode being used frequently by this user.



3.3.4 Evidence of malicious activity

During our analysis of the LiveJournal dataset using *SocialSEAL*, we were able to find some evidence of malicious activity. We examined the posts made by the user **eas7**, who was found to have the highest Risk and Influence scores among all users in the dataset. We discovered that this user has authored many articles on cybersecurity and seems to be very knowledgeable. Some of the findings through our investigation on **eas7** include -

- This user has shared an article on exploiting the Key Frame Buffer Overflow in IE6. The vulnerability which this exploit can be used for is listed under CVE-2006-4446.

- The user has shared information on tools which could be used to hack Bluetooth devices such as BlueBugger, BTcrack etc.
- The user seems to be closely associated with the anonymous Internet forum project <http://geeklife.ru>, which was active between the years 2007 and 2011. On examining this website using the Internet Archive (<http://web.archive.org>), we found that users on this website could anonymously share information on hacking and possibly engage in other malicious activities. The website has several mirrors such as <http://geek.grayhat.ru> and <http://geek.rootkits.ru>. There have also been some reports of spam originating from some users of this website.
- The user has shared information on a vulnerability in Yahoo messenger listed under CVE-2007-4901 and shared links to an exploit, which can be accessed through the Internet Archive:
 - http://shinnai.altervista.org/exploits/txt/TXT_KJDPaI2IIM5P9PP6N6dl.html
 - <http://milw0rm.com/exploits/4428>
- The user owns the personal website <http://eas7.ru> which provides several tutorials for authoring exploits.

Similarly, we found another user named **cr4sh** who seems to be a close relationship with **eas7** and appears to be a very popular hacker. This user has a very active online presence and is the author of a few blogs which talk about the exploits and rootkits that he/she has authored. We found a post authored by **cr4sh** that talks about a spyware named Agent.btz [15] which affects the Windows OS. Another user named **bhcrew** has the most number of friends in the dataset and is the person who maintains the hacking community which goes by the name **bh_crew**. This user has authored many online hacking magazines and has distributed those magazines through the community website <http://bhcrew.org>, which is no longer accessible but can be viewed through the Internet Archive. One of the friends in **bhcrew** who goes by the username **t1m** has authored an article that contains a link on how to create viruses, which is hosted on the once-popular virus exchange underground forum <http://vx.netlux.org>.

Upon examining all the links shared by the various users in the dataset, we were able to find some links which were identified as Malware by the Google Safe Browsing API. We also found some posts talking about Phrack [16], a popular online hacking magazine and we examined underground websites such as <http://mazafaka.ru>, which some users in the dataset are members of. Consequently, we were able to find some evidence of malicious activity in the dataset analyzed using *SocialSEAL*. Through deeper investigation, we believe it is quite possible that more actionable evidence can be unearthed in the next phase of this project.

4 Conclusion and Future Work

In this report of the seed project, we have presented the *SocialSEAL* framework and tool, which could be used by security analysts in fighting cyber threats. The proof-of-concept tool demonstrated how our framework could use of social network analysis techniques to identify the influential players among the online underground hacking community. Our evaluation of *SocialSEAL* on real-world data has also shown the feasibility and effectiveness of *SocialSEAL*

and demonstrated its capabilities for end-users to analyze the data with more meaningful and actionable insight.

In the next phase of this project, we plan to make *SocialSEAL* more scalable and provide the ability to track multiple OSNs so that it can gradually reduce the manual effort involved in analyzing large datasets. In addition, we will further improve and refine our framework and tool to handle and analyze new data in a more dynamic and real-time manner.

Also, we will pursue to collect diverse data related to cyber threats. For example, we are currently collecting twelve forums between December 2005 and July 2011, that cover common topics in stolen data markets, including 'carding', 'dump', 'purchase', 'sale', and 'cvv'. We will explore the contents of one Russian Speaking Carder subforums and continuously collect Russian language forums via links shared by users and the marketplace data in this time frame. We believe it will show its unique value, because the earliest documented online cybercrime markets emerged around 2003 and 2011 is the year Silk Road was launched and Bitcoin gained its popularity. Also, we plan to pull the main jihadi forum from the Ansar Al-Mujahideen English Forum (AMEF). Interestingly, this forum includes various discussions that deal with cyber hacking tools and mechanisms.

Immediate support for the next phase of this project is needed since the preliminary version of proof-of-concept prototype has been implemented and it is expected that the next phase of this research project will greatly enhance the progress achieved from the current seed project.

5 References

1. P. Bacher, T. Holz, M. Kotter, and G. Wicherski, "Know your Enemy: Tracking Botnets—Using honeynets to learn more about Bots," 2005.
2. T.J. Holt, J.W. Burruss and A.M. Bossler, "Social Learning and Cyber-Deviance: Examining the Importance of a Full Social Learning Model in the Virtual World," *Journal of Crime and Justice*, p. 33, 2010.
3. E. Athanasopoulos, A. Makridakis, S. Antonatos, D. Antoniadis, S. Ioannidis, K. Anagnostakis, and E. Markatos, "Antisocial Networks: Turning a Social Network into a Botnet," In *Proc. of the 11th International Conference on Information Security (ISC)*, Springer, 2008.
4. R. Lemos, "Script kiddies: The Net's cybergangs," Accessed at: <http://www.zdnet.com/script-kiddies-the-nets-cybergangs-3002080125/>, 2000.
5. Z. Zhao, G.-J. Ahn, H. Hu and D. Mahi, "SocialImpact: Systematic Analysis of Underground Social Dynamics," In *Proc. of 17th European Symposium on Research in Computer Security (ESORICS)*, 2012.
6. N. Agarwal, H. Liu, L. Tang, and P. Yu, "Identifying the influential bloggers in a community," In *Proc. of the 1st International Conference on Web Search and Web Data Mining (WSDM)*. ACM, 2008.
7. A. Singhal, "Modern Information Retrieval: A Brief Overview," *Bulletin of the IEEE Computer Society Technical Committee on Data Engineering* 24 (4): 35–43. 2001.
8. C. Vicknair, M. Macias, Z. Zhao, X. Nan, Y. Chen and D. Wilkins, "A comparison of a graph database and a relational database: a data provenance perspective," In *Proc. of the 48th Annual Southeast Regional Conference, ACM SE*, 2010.

9. S. Maple, "The Curious Coder's Java Web Frameworks Comparison: Spring MVC, Grails, Vaadin, GWT, Wicket, Play, Struts and JSF," Accessed at: <http://zeroturnaround.com/rebellabs/the-curious-coders-java-web-frameworks-comparison-spring-mvc-grails-vaadin-gwt-wicket-play-struts-and-jsf/>, 2013.
10. M. J. Bannister, D. Eppstein, M. T. Goodrich and L. Trott, "Force-directed graph drawing using social gravity and scaling," Proc. 20th International Symposium on Graph Drawing, 2012.
11. M. Halvey and M. T. Keane, "An Assessment of Tag Presentation Techniques," In Proc. of 16th International World Wide Web Conference, 2007.
12. L. Barrett, "Russia, Brazil Lead Cyber Attack Barrage," <http://www.esecurityplanet.com/>, 2010.
13. "Glossary of Security Terms," Accessed at: <http://www.sans.org/security-resources/glossary-of-terms/>
14. "A Glossary of Common Cybersecurity Terminology," Accessed at: <http://niccs.us-cert.gov/glossary>
15. Spyware Report, Accessed at http://www.avira.com/en/support-threats-description/tid/2820/tr_dldr.agent.awf.14.html
16. T. King, "Phrack Introduction," Accessed at: <http://phrack.org/issues/1/1.html#article>